

Please type a plus sign (+) inside this box → ☒

Approved for use through 09/30/00. OMB 0651-0032  
Patent and Trademark Office: U.S. DEPARTMENT OF COMMERCE

Under the Paperwork Reduction Act of 1995, no persons are required to respond to a collection of information unless it displays a valid OMB control number

# UTILITY PATENT APPLICATION TRANSMITTAL

(Only for new nonprovisional applications under 37 CFR 1.53(b))

Attorney Docket No.

TI-25489.1

First Named Inventor or Application Identifier

Yifan Gong

Title

Source Normalization Training for HMM Modeling of Speech

Express Mail Label No.

EL547745323US

## APPLICATION ELEMENTS

See MPEP Chapter 600 concerning utility patent application contents

## ADDRESS TO:

Assistant Commissioner for Patents  
Box Patent Application  
Washington, DC 20231

1. ☒ \*Fee Transmittal Form (e.g., PTO/SB/17)  
(Submit an original, and a duplicate for fee processing)
2. ☒ Specification [Total Pages **21**]  
(preferred arrangement set forth below)  
- Descriptive title of the invention  
- Cross References to Related Applications  
- Statement Regarding Fed sponsored R&D  
- Reference to Microfiche Appendix  
- Background of the Invention  
- Brief Summary of the Invention  
- Brief Description of the Drawings (if filed)  
- Detailed Description  
- Claim(s)  
- Abstract of the Disclosure
3. ☒ Drawing(s) (35 USC d113) [Total Sheets **3**]
4. Oath or Declaration [Total Pages **1**]  
a. ☐ Newly Executed (original or copy)  
b. ☒ Copy from a prior application (37 CFR §1.63(d))  
(for continuation/divisional with Box 17 completed)  
[Note Box 5 below]  
i. ☐ **DELETION OF INVENTOR(S)**  
Signed statement attached deleting inventor(s)  
named in the prior application,  
see 37 CFR §1.63(d)(2) and 1.33(b).
5. ☒ Incorporation By Reference (useable if Box 4b is checked)  
The entire disclosure of the prior application, from which a copy of  
the oath or declaration is supplied under Box 4b, is considered as  
being part of the disclosure of the accompanying application and is  
hereby incorporated by reference therein.

6. ☐ Microfiche Computer Program (Appendix)
7. Nucleotide and/or Amino Acid Sequence Submission  
(if applicable, all necessary)  
a. ☐ Computer Readable Copy  
b. ☐ Paper Copy (identical to computer copy)  
c. ☐ Statement verifying identical of above copies

## ACCOMPANYING APPLICATION PARTS

8. ☐ Assignment Papers (cover sheet & Documents(s))
9. ☐ 37 CFR §3.73(b) Statement (when there is an assignee) ☒ Power of Attorney
10. ☐ English Translation Document (if applicable)
11. ☐ Information Disclosure Statement (IDS)/PTO-1449 ☐ Copies of IDS Citations
12. ☒ Preliminary Amendment
13. ☒ Return Receipt Postcard (MPEP 503)  
(Should be specifically itemized)
14. ☐ \*Small Entity Statement(s) ☐ Statement filed in prior application  
(PTO/SB/09-12) Status still proper and desired
15. ☐ Certified Copy of Priority Document(s)  
if foreign priority is claimed)
16. ☐ Other.

\*A new statement is required to be entitled to pay small entity fees, except  
where one has been filed in a prior application and is being relied upon

17. If a **CONTINUING APPLICATION**, check appropriate box and supply the requisite information below and in a preliminary amendment:

☐ Continuation

☒ Divisional

☐ Continuation-in-part (CIP)

of prior application No: 09/134,775

Prior application information: Examiner A. Azad

Group / Art Unit: 2741

## 18. CORRESPONDENCE ADDRESS



Customer Number or Bar Code Label

(Insert Customer No. or Attach bar code label  
here)



Correspondence address below

NAME	Robert L. Troike		
ADDRESS	P.O. Box 655474 MS 3999		
CITY	Dallas	STATE	TX
COUNTRY	US	TELEPHONE	202-639-7710
		ZIP CODE	75265-5474
		FAX	202-639-7890

Name (Print/Type) Robert L. Troike Registration No. (Attorney/Agent) 24,183

Signature

*Robert L. Troike*

Date

6/3/00

**Burden Hour Statement:** This form is estimated to take 0.2 hours to complete. Time will vary depending upon the needs of the individual case. Any comments on the amount of time you are required to complete this form should be sent to the Chief Information Officer, Patent and Trademark Office, Washington, DC 20231. DO NOT SEND FEES OR COMPLETED FORMS TO THIS ADDRESS. SEND TO: Assistant Commissioner for Patents, Box Patent Application, Washington, DC 20231.

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re the Application of:

Yifan Gong

Serial No: TBD

Filed: Herewith

For: Source Normalization Training for HMM Modeling of Speech

TI-25489.1

Examiner: TBD

Art Unit: TBD

PRELIMINARY AMENDMENT

Assistant Commissioner for Patents

Washington, D. C. 20231

Dear Sir:

"EXPRESS MAILING" Mailing Label No. EL547745323US I hereby certify that this paper is being deposited with the U.S. Postal Service Express Mail Post Office to Addressee Service under 37 CFR 1.10 on the date shown below and is addressed to the Assistant Commissioner for Patents, Washington, D.C. 20231.

*Robert L. Troike*

Robert L. Troike, Reg. No. 24, 183

*6/7/00*

Date

Prior to the examination of the above-identified application, please amend as follows:

IN THE SPECIFICATION:

Page 1, insert before the first sentence:

--This application is a divisional of prior application number 09/134,775, filed 08/15/98--

**IN THE CLAIMS:**

Please cancel claims 1-5, amend claim 6 and add claims 7 and 8 as follows:

6. (Amended) An improved speech recognition system comprising:

a speech recognizer; and

a source normalization model derived by application of an estimation maximization algorithm with explicit separation of source information and environment distortion factors in an unsupervised manner

7. A recognition system comprising:

an input signal;

a stored reference;

and a comparator for comparing said input signal with said stored reference;

said stored reference has at least two components, one representing signal source and the other representing transformations for a number of environments;

said reference is identified by source normalization training, which consists in performing the following steps of:

- (a) determining a new set of signal source representation or at least part of the representation that reduces the distance between the new reference and a training signal, given training signals and current transformations and
- (b) for each environment, determine a new transformation or at least part of it that, jointly with the signal source representation, reduces the distance between the new reference and the training signal where said environment represents either a label associated with training signal or a class of distortion sources.

[illegible]

receiving an input signal; and

Respectfully submitted,

Robert L. Troike  
Attorney for Applicant  
Reg. No. 24,183

Texas Instruments Incorporated  
P. O. Box 655474, MS 3999  
Dallas, TX 75265  
(202) 639-7710  
Fax: (202) 639-7890

# SOURCE NORMALIZATION TRAINING FOR HMM MODELING OF SPEECH

## TECHNICAL FIELD OF THE INVENTION

This invention relates to training for HMM modeling of speech and more particularly to removing environmental factors from speech signal during the training procedure.

## BACKGROUND OF THE INVENTION

In the present application we refer to environment as speaker, handset or microphone, transmission channel, noise background conditions, or combination of these as the environment. A speech signal can only be measured in a particular environment. Speech recognizers suffer from environment variability because trained model distributions may be biased from testing signal distributions because environment mismatch and trained model distributions are flat because they are averaged over different environments.

The first problem, the environmental mismatch, can be reduced through model adaptation, based on some utterances collected in the testing environment. To solve the second problem, the environmental factors should be removed from the speech signal during the training procedure, mainly by source normalization.

In the direction of source normalization, speaker adaptive training uses linear regression (LR) solutions to decrease inter-speaker variability. See for example, T. Anastasakos, et al. entitled, "A compact model for speaker-adaptive training," *International Conference on Spoken Language Processing*, Vol. 2, October 1996. Another technique models mean-vectors as the sum of a speaker-independent bias and a speaker-dependent vector. This is found in A. Acero, et al. entitled, "Speaker and Gender Normalization for Continuous-Density Hidden Markov Models," in *Proc. Of IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 342-345, Atlanta, 1996. Both of these techniques require explicit label of the classes. For example, speaker or gender of the utterance during the training. Therefore, they can not be used to train clusters of classes, which represent acoustically close speaker, hand set or microphone, or background noises. Such inability of discovering clusters may be a disadvantage in application.

## SUMMARY OF THE INVENTION

In accordance with one embodiment of the present invention, we provide a maximum likelihood (ML) linear regression (LR) solution to the environment  
5 normalization problem, where the environment is modeled as a hidden (non-observable) variable. An EM-Based training algorithm can generate optimal clusters of environments and therefore it is not necessary to label a database in terms of environment. For special cases, the technique is compared to utterance-by-utterance cepstral mean normalization (CMN) technique and show performance improvement on a noisy speech telephone  
10 database.

In accordance with one embodiment of the present invention under maximum-likelihood (ML) criterion, by application of EM algorithm and extension of Baum-Welch forward and backward variables and algorithm, we obtained joint solution to the parameters  
15 for the source normalization, i.e. the canonical distributions, the transformations and the biases.

These and other features of the invention that will be apparent to those skilled in the art from the following detailed description of the invention, taken together with the  
20 accompanying drawings.

## DESCRIPTION OF THE DRAWINGS

Fig. 1 is a block diagram of the system according to one embodiment of the present invention;

5 Fig. 2 illustrates a speech model;

Fig. 3 illustrates a Gaussian distribution;

Fig. 4 illustrates distortions in the distribution caused by different environments;

Fig. 5 is a more detailed flow diagram of the process according to one embodiment of the present invention; and

10 Fig. 6 is a recognizer according to an embodiment of the present invention using a source normlization model.



## DESCRIPTION OF PREFERRED EMBODIMENTS OF THE PRESENT INVENTION

The training is done on a computer workstation which is illustrated in Fig. 1  
5 having a monitor 11, a computer workstation 13, a keyboard 15, and a mouse or other  
interactive device 15a as shown in Fig. 1. The system maybe connected to a separate  
database represented by database 17 in Fig. 1 for storage and retrieval of models.

By the term "training" we mean herein to fix the parameters of the speech models  
10 according to an optimum criterion. In this particular case, we use HMM (Hidden Markov  
Models) models. These models are as represented in Fig. 2 with states A, B, and C and  
transitions E, F, G, H, I and J between states. Each of these states has a mixture of  
Gaussian distributions 18 represented by Fig. 3. We are training these models to account  
for different environments. By environment we mean different speaker, handset,  
15 transmission channel, and noise background conditions. Speech recognizers suffer from  
environment variability because trained model distributions may be biased from testing  
signal distributions because of environment mismatch and trained model distributions are  
flat because they are averaged over different environments. For the first problem, the  
environmental mismatch can be reduced through model adaptation, based on utterances  
20 collected in the testing environment. Applicant's teaching herein is to solve the second  
problem by removing the environmental factors from the speech signal during the  
training procedure. This is source normalization training according to the present  
invention. A maximum likelihood (ML) linear regression (LR) solution to the

environmental problem is provided herein where the environment is modeled as hidden-  
(non observable) variable.

A clean speech pattern distribution 40 will undergo complex distortion with  
different environments as shown in Fig. 4. The two axes represent two parameters which  
may be, for example, frequency, energy, formant, spectral, or cepstral components. The  
Fig. 4 illustrates a change at 41 in the distribution due to background noise or a change in  
speakers. The purpose of the application is to model the distortion.

The present model assumes the following: 1) the speech signal  $x$  is generated by  
Continuous Density Hidden Markov Model (CDHMM), called source distributions; 2)  
before being observed, the signal has undergone an environmental transformation, drawn  
from a set of transformations, where  $W_{je}$  be the transformation on the HMM state  $j$  of the  
environment  $e$ ; 3) such a transformation is linear, and is independent of the mixture  
components of the source; and 4) there is a bias vector  $b_{ke}$  at the  $k$ -th mixture component  
due to environment  $e$ .

What we observe at time  $t$  is:

$$o_t = W_{je}x_t + b_{ke} \quad (1)$$

Our problem now is to find, in the maximum likelihood (ML) sense, the optimal  
source distributions, the transformation and the bias set.

In the prior art (A. Acero, et al. cited above and T. Anastasakos, et al. cited above), the environment  $e$  must be explicit, e.g.: speaker identity, male/female. This work overcomes this limitation by allowing an arbitrary number of environments which are optimally trained.

Let  $N$  be the number of HMM states,  $M$  be the mixture number,  $L$  be the number of environments,  $\Omega_s \triangleq \{1, 2, \dots, N\}$  be the set of states  $\Omega_m \triangleq \{1, 2, \dots, M\}$  be the set of mixture indicators, and  $\Omega_e \triangleq \{1, 2, \dots, L\}$  be the set of environmental indicators.

For an observed speech sequence of  $T$  vectors:  $O \triangleq o_1^T \triangleq (o_1, o_2, \dots, o_T)$ , we introduce state sequence  $\Theta \triangleq \{\theta_1, \dots, \theta_T\}$  where  $\theta_t \in \Omega_s$ , mixture indicator sequence  $\Xi \triangleq (\xi_1, \dots, \xi_T)$  where  $\xi_t \in \Omega_m$ , and environment indicator sequence  $\Phi \triangleq (\phi_1, \dots, \phi_T)$  where  $\phi_t \in \Omega_e$ . They are all unobservable. Under some additional assumptions, the joint probability of  $O, \Theta, \Xi$ , and  $\Phi$  given model  $\lambda$  can be written as:

$$p(O, \Theta, \Xi, \Phi | \lambda) = u_{\theta_1} \prod_{t=1}^T c_{\theta_t \xi_t} b_{\theta_t \xi_t \phi_t}(o_t) a_{\theta_t \theta_{t+1}} l_{\phi_t} \quad (2)$$

where

$$b_{jke}(o_t) \triangleq p(o_t | \theta_t = j, \xi_t = k, \phi_t = e, \lambda) \quad (3)$$

$$= N(o_i; W_{je} \mu_{jk} + b_{ke}, \sum_{jk}), \quad (4)$$

$$u_i \triangleq p(\theta_1 = i), \quad a_{ij} \triangleq p(\theta_{t+1} = j | \theta_t = i) \quad (5)$$

$$c_{jk} \triangleq p(\xi_t = k | \theta_t = j, \lambda), \quad l_e \triangleq p(\varphi = e | \lambda) \quad (6)$$

5

Referring to Fig. 1, the workstation 13 including a processor contains a program as illustrated that starts with an initial standard HMM model 21 which is to be refined by estimation procedures using Baum-Welch or Estimation-Maximization procedures 23 to get new models 25. The program gets training data at database 19 under different environments and this is used in an iterative process to get optimal parameters. From this model we get another model 25 that takes into account environment changes. The quantities are defined by probabilities of observing a particular input vector at some particular state for a particular environment given the model.

10

15

The model parameters can be determined by applying generalized EM-procedure with three types of hidden variables: state sequence, mixture component indicators, and environment indicators. (A. P. Dempster, N. M. Laird, and D. B. Rubin, entitled "Maximum Likelihood from Incomplete Data via the EM Algorithm," *Journal of the Royal Statistical Society*, 39 (1): 1-38, 1977.) For this purpose, Applicant teaches the CDHMM formulation from B. Juang, "Maximum-Likelihood Estimation for Mixture Multivariate Stochastic Observation of Markov Chains" (*The Bell System Technical*

20

*Journal*, pages 1235-1248, July-August 1985) to be extended to result in the following paragraphs: Denote:

$$\alpha_t(j, e) \triangleq p(\alpha_1^t, \theta_t = j, \varphi = e | \bar{\lambda}) \quad (7)$$

$$\beta_t(j, e) \triangleq p(\alpha_{t+1}^T | \theta_t = j, \varphi = e, \bar{\lambda}) \quad (8)$$

$$\gamma_t(j, k, e) \triangleq p(\theta_t = j, \xi_t = k, \varphi = e | O, \bar{\lambda}) \quad (9)$$

The speech is observed as a sequence of frames (a vector). Equations 7, 8, and 9 are estimations of intermediate quantities. For example, in equation 7 is the joint probability of observing the frames from times 1 to t at the state j at time t and for the environment of e given the model  $\lambda$ .

The following re-estimation equations can be derived from equations 2, 7, 8, and 9.

For the EM procedure 23, equations 10-21 are solutions for the quantities in the model.

*Initial state probability:*

$$u_i = \frac{1}{R} \sum_{r=1}^R \frac{\sum_{e \in \Omega_e} \alpha_1^r(i, e) \beta_1^r(i, e)}{\sum_{i \in \Omega_s} \sum_{e \in \Omega_e} \alpha_1^r(i, e) \beta_1^r(i, e)} \quad (10)$$

with R the number of training tokens.

Transition probability:

$$a_{ij} = \frac{\bar{a}_{ij} \sum_{r=1}^R \frac{1}{p(O^r|\bar{\lambda})} \sum_{e \in \Omega_e} \sum_{t=1}^T \alpha_t^r(i, e) b_{j,e}(o_{t+1}^r) \beta_{t+1}^r(j, e)}{\sum_{r=1}^R \frac{1}{p(O^r|\bar{\lambda})} \sum_{e \in \Omega_e} \sum_{t=1}^T \alpha_t^r(i, e) \beta_t^r(i, e)} \quad (11)$$

Mixture Component probability: (Mixture probability is where there is a mixture of Gaussian distributions)

$$c_{jk} = \frac{\sum_{r=1}^R \sum_{e \in \Omega_e} \sum_{t=1}^T \gamma_t^r(j, k, e)}{\sum_{r=1}^R \frac{1}{p(O^r|\bar{\lambda})} \sum_{e \in \Omega_e} \sum_{t=1}^T \alpha_t^r(j, e) \beta_t^r(j, e)} \quad (12)$$

Environment probability:

$$l_e = \frac{1}{R} \sum_{r=1}^R \frac{\sum_{j \in \Omega_j} \alpha_T^r(j, e)}{\sum_{e \in \Omega_e} \sum_{j \in \Omega_j} \alpha_T^r(j, e)} \quad (13)$$

Mean vector and bias vector: We introduce:

$$\rho(j, k, e) \triangleq \sum_{r=1}^R \sum_{t=1}^T \gamma_t^r(j, k, e) o_t^r \quad (14)$$

$$g(j, k, e) \triangleq \sum_{r=1}^R \sum_{t=1}^T \gamma_t^r(j, k, e) \quad (15)$$

and

$$G_{ke} = \sum_{j \in \Omega_j} g(j, k, e) \sum_{jk}^{-1} \quad (16)$$

$$E_{jke} = g(j, k, e) W_{je}' \sum_{jk}^{-1} \quad (17)$$

$$F_{jk} = \sum_{e \in \Omega_e} E_{jke} W_{je} \quad (18)$$

$$a_{jk} = \sum_{e \in \Omega_e} W_{je}' \sum_{jk}^{-1} \rho(j, k, e) \quad (19)$$

$$c_{ke} = \sum_{j \in \Omega_s} \sum_{jk}^{-1} \rho(j, k, e). \quad (20)$$

5 Assuming  $W_{je} = \overline{W_{je}}$  and  $\sum_{jk}^{-1} = \overline{\sum_{jk}^{-1}}$ , for a given k, we have N+L equations:

$$\sum_{e \in \Omega_e} E_{jke} b_{ke} + F_{jk} \mu_{jk} = a_{jk} \quad \forall j \in \Omega_s \quad (21)$$

$$G_{ke} b_{ke} + \sum_{j \in \Omega_s} H_{jke} \mu_{jk} = c_{ke} \quad \forall e \in \Omega_e \quad (22)$$

10 These equations 21 and 22 are solved jointly for mean vectors and bias vectors.

Therefore  $\mu_{jk}$  and  $b_{ke}$  can be simultaneously obtained by solving the linear system of N+L variables.

15 *Covariance:*

$$\sum_{jk} = \frac{\sum_{e \in \Omega_e} \sum_{r=1}^R \sum_{t=1}^{T'} \gamma_t'(j, k, e) \delta_t'(j, k, e) \delta_t'(j, e, k)'}{\sum_{e \in \Omega_e} g(j, k, e)} \quad (23)$$

where  $\delta_i^r(j, k, e) \triangleq o_i^r - W_{je} \mu_{jk} - b_{ke}$

*Transformation:* We assume covariance matrix to be diagonal:  $\sum_{jk}^{-1(m,n)} = 0$  if  $n \neq m$ .

For the line  $m$  of transformation  $W_{je}$ , we can derive (see for example C. J. Leggetter, et al., entitled "Maximum Likelihood Linear Regression for Speaker Adaptation of Continuous Density HMMs" *Computer, Speech and Language*, 9(2): 171-185, 1995.):

$$Z_{je}^{(m)} = W_{je}^{(m)} R_{je}^{(m)} \quad (24)$$

which is a linear system of  $D$  equations, where:

$$Z_{je}^{(m,n)} \triangleq \sum_{k \in \Omega_m} \sum_{jk}^{-1(m,m)} \mu_{jk}^{(n)} \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j, k, e) (o_i^r - b_{ke})^{(m)} \quad (25)$$

$$R_{je}^{(p,n)}(m) \triangleq \sum_{k \in \Omega_m} \sum_{jk}^{-1(m,m)} \mu_{jk}^{(p)} \mu_{jk}^{(n)} \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j, k, e). \quad (26)$$

Assume the means of the source distributions ( $\mu_{jk}$ ) are constant, then the above set of source normalization formulas can also be used for model adaptation.

The model is specified by the parameters. The new model is specified by the new parameters.



As illustrated in Figs. 1 and 5, we start with an initial as standard model 21 such as the CDHMM model with initial values. This next step is the Estimation Maximization 23 procedure starting with (Step 23a) equations 7-9 and re-estimation (Step 23b) equations 10-13 for initial state probability, transition probability, mixture component probability and environment probability.

The next step (23c) to derive means vector and bias vector by introducing two additional equations 14 and 15 and equation 16-20. The next step 23a is to apply linear equations 21 and 22 and solve 21 and 22 jointly for mean vectors and bias vectors and at the same time calculate the variance using equation 23. Using equation 24 which is a system of linear equations will solve for transformation parameters using quantities given by equation 25 and 26. Then we have solved for all the model parameters. Then one replaces the old model parameters by the newly calculated ones (Step 24). Then the process is repeated for all the frames. When this is done for all the frames of the database a new model is formed and then the new models are re-evaluated using the same equation until there is no change beyond a predetermined threshold (Step 27).

After a source normalization training model is formed, this model is used in a recognizer as shown in Fig. 6 where input speech is applied to a recognizer 60 which used the source normalized HMM model 61 created by the above training to achieve the response.

The recognition task has 53 commands of 1-4 words. ("call return", "cancel call return", "selective call forwarding", etc.). Utterances are recorded through telephone lines, with a diversity of microphones, including carbon, electret and cordless microphones and hands-free speaker-phones. Some of the training utterances do not correspond to their transcriptions. For example: "call screen" (cancel call screen), "matic call back" (automatic call back), "call tra" (call tracking).

The speech is 8kHz sampled with 20ms frame rate. The observation vectors are composed of LPCC (Linear Prediction Coding Coefficients) derived 13-MFCC (Mel-Scale Cepstral Coefficients) plus regression based delta MFCC. CMN is performed at the utterance level. There are 3505 utterances for training and 720 for speaker-independent testing. The number of utterances per call ranges between 5-30.

Because of data sparseness, besides transformation sharing among states and mixtures, the transformations need to be shared by a group of phonetically similar phones. The grouping, based on an hierarchical clustering of phones, is dependent on the amount of training (SN) or adaptation (AD) data, i.e., the larger the number of tokens is, the larger the number of transformations. Recognition experiments are run on several system configurations:

- BASELINE applies CMN utterance-by-utterance. This simple technique will remove channel and some long term speaker specificities, if the duration of the utterance is long enough, but can not deal with time domain additive noises.

- SN performs source-normalized HMM training, where the utterances of a phone-call are assumed to have been generated by a call-dependent acoustic source. Speaker, channel and background noise that are specific to the call is then removed by MLLR. An HMM recognizer is then applied using source parameters. We evaluated a special case, where each call is modeled by one environment.

- AD adapts traditional HMM parameters by unsupervised MLLR. 1. Using current HMMs and task grammar to phonetically recognize the test utterances, 2. Mapping the phone labels to a small number (N) of classes, which depends on the amount of data in the test utterances, 3. Estimating the LR using the N-classes and associated test data, 4. Recognizing the test utterances with transformed HMM. A similar procedure has been introduced in C. J. Legetter and P. C. Woodland. "Maximum likelihood linear regression for speaker adaptation of continuous density HMMs." *Computer, Speech and Language*, 9(2):171-185, 1995.

- SN+AD refers to AD with initial models trained by SN technique.

Based on the results summarized in Table 1, we point out:

- For numbers of mixture components per state smaller than 16, SN, AD, and SN+AD all give consistent improvement over the baseline configuration.
- For numbers of mixture components per state smaller than 16, SN gives about 10% error reduction over the baseline. As SN is a training procedure which does not require any change to the recognizer, this error reduction mechanism immediately benefits applications.
- For all tested configurations, AD using acoustic models trained with SN procedure always gives additional error reduction.
- The most efficient case of SN+AD is with 32 components per state, which reduces error rate by 23%, resulting 4.64% WER on the task.

	4	8	16	32
baseline	7.85	6.94	6.83	5.98
SN	7.53	6.35	6.51	6.03
AD	7.15	6.41	5.61	5.87
SN+AD	6.99	6.03	5.41	4.64

Table 1: Word error rate (%) as function of test configuration and number of mixture components per state.

Although the present invention and its advantages have been described in detail, it should be understood that various changes, substitutions and alterations can be made herein

without departing from the spirit and scope of the invention as defined by the appended claims.

TI-25489 (Page 17)

**WHAT IS CLAIMED IS:**

1. A method of source normalization training for HMM modeling of speech comprising the steps of:

5 (a) providing an initial model;

(b) on said initial model or following new models performing the following steps to get a new model:

b<sub>1</sub>) estimation of intermediate quantities;

b<sub>2</sub>) performing re-estimation to determine initial state probability, transition probability, mixture component probability and environment probability;

b<sub>3</sub>) deriving mean vector and bias vector;

b<sub>4</sub>) solving jointly for mean vector and bias vector using linear equations and determining variances and transformation; and

b<sub>5</sub>) replacing old model parameters for the calculated ones; and

c) determining after a new model is formed if it differs significantly from the previous model and if so repeating steps b<sub>1</sub> - b<sub>5</sub>.

2. The method of Claim 1 wherein in step b<sub>1</sub> estimation intermediate quantities is determined by  $\alpha_i(j, e) \triangleq p(\alpha'_i, \theta_i = j, \varphi = e | \bar{\lambda})$ ,

$\beta_i(j, e) \triangleq p(\alpha_{i+1}^T | \theta_i = j, \varphi = e, \bar{\lambda})$ , and  $\gamma_i(j, k, e) \triangleq p(\theta_i = j, \xi_k = k, \varphi = e | O, \bar{\lambda})$ .

3. The method of Claim 2 wherein step b<sub>2</sub> the initial state probability is determined by  $u_i = \frac{1}{R} \sum_{r=1}^R \frac{\sum_{e \in \Omega_e} \alpha_i^r(i, e) \beta_i^r(i, e)}{\sum_{i \in \Omega_s} \sum_{e \in \Omega_e} \alpha_i^r(i, e) \beta_i^r(i, e)}$ , transition probability is determined by

$$a_{ij} = \frac{\bar{a}_{ij} \sum_{r=1}^R \frac{1}{p(O^r|\bar{\lambda})} \sum_{e \in \Omega_e} \sum_{t=1}^{T^r} \alpha_t^r(i, e) b_{j,e}(o_{t+1}^r) \beta_{t+1}^r(j, e)}{\sum_{r=1}^R \frac{1}{p(O^r|\bar{\lambda})} \sum_{e \in \Omega_e} \sum_{t=1}^{T^r} \alpha_t^r(i, e) \beta_t^r(i, e)}, \text{ mixture component probability is}$$

determined by  $c_{jk} = \frac{\sum_{r=1}^R \sum_{e \in \Omega_e} \sum_{t=1}^{T^r} \gamma_t^r(j, k, e)}{\sum_{r=1}^R \frac{1}{p(O^r|\bar{\lambda})} \sum_{e \in \Omega_e} \sum_{t=1}^{T^r} \alpha_t^r(j, e) \beta_t^r(j, e)}$ , and environment

probability is determined by  $l_e = \frac{1}{R} \sum_{r=1}^R \frac{\sum_{j \in \Omega_s} \alpha_r^r(j, e)}{\sum_{e \in \Omega_e} \sum_{j \in \Omega_s} \alpha_r^r(j, e)}$ .

4. The method of Claim 2 wherein step b<sub>3</sub> deriving mean vector and bias vector is determined by  $\rho(j, k, e) \triangleq \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j, k, e) o_t^r$ ;  $g(j, k, e) \triangleq \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j, k, e)$ ,

$$G_{ke} = \sum_{j \in \Omega_s} g(j, k, e) \sum_{jk}^{-1}, \quad E_{jke} = g(j, k, e) W_{je}' \sum_{jk}^{-1}, \quad F_{jk} = \sum_{e \in \Omega_e} E_{jke} W_{je},$$

$$a_{jk} = \sum_{e \in \Omega_e} W_{je}' \sum_{jk}^{-1} \rho(j, k, e), \text{ and } c_{ke} = \sum_{j \in \Omega_s} \sum_{jk}^{-1} \rho(j, k, e).$$

5. The method of Claim 2 wherein step b<sub>4</sub> equations  $\sum_{e \in \Omega_e} E_{jke} b_{ke} + F_{jk} \mu_{jk} = a_{jk} \quad \forall j \in \Omega_s$  and  $G_{ke} b_{ke} + \sum_{j \in \Omega_s} H_{jke} \mu_{jk} = c_{ke} \quad \forall e \in \Omega_e$  are used

for solving jointly and equation  $\sum_{jk} = \frac{\sum_{e \in \Omega_e} \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j, k, e) \delta_t^r(j, k, e) \delta_t^r(j, e, k)'}{\sum_{e \in \Omega_e} g(j, k, e)}$

is used to determine variance and equations  $Z_{je}^{(m)} = W_{je}^{(m)} R_{je}(m)$ ,

$$Z_{je}^{(m,n)} \triangleq \sum_{k \in \Omega_m} \sum_{jk}^{-1(m,m)} \mu_{jk}^{(n)} \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j, k, e) (o_t^r - b_{ke})^{(m)}, \quad \text{and}$$

$$R_{je}^{(p,n)}(m) \triangleq \sum_{k \in \Omega_m} \sum_{jk}^{-1(m,m)} \mu_{jk}^{(p)} \mu_{jk}^{(n)} \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j, k, e). \text{ are used to determine transformation.}$$

- 5
6. An improved speech recognition system comprising:
- a speech recognizer; and
- a source normalization model derived by application of an estimation maximization algorithm.



## ABSTRACT OF THE DISCLOSURE

A maximum likelihood (ML) linear regression (LR) solution to environment normalization is provided where the environment is modeled as a hidden (non-observable) variable. By application of an expectation maximization algorithm and extension of Baum-Welch forward and backward variables (Steps 23a-23d) a source normalization is achieved such that it is not necessary to label a database in terms of environment such as speaker identity, channel, microphone and noise type.

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Application of

Yifan Gong

Serial No.: TBD

Filed: Herewith

For: Source Normalization Training for HMM Modeling of Speech

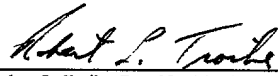
TI-25489.1

Examiner: TBD

Art Unit: TBD

LETTER TO THE OFFICIAL DRAFTSPERSON

"EXPRESS MAILING" Mailing Label No. EL547745323US I hereby certify that this paper is being deposited with the U.S. Postal Service Express Mail Post Office to Addressee Service under 37 CFR 1.10 on the date shown below and is addressed to the Assistant Commissioner for Patents, Washington, D.C. 20231.

 6/7/00  
Robert L. Troike, Reg. No. 24, 183 Date

Assistant Commissioner for Patents

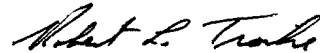
Washington, D.C. 20231

Sir:

Enclosed are **THREE (3)** sheets of formal drawings for the above-referenced case. Please charge any necessary fees to Deposit Account No. 20-0668 of Texas Instruments Incorporated.

This sheet is enclosed in triplicate.

Respectfully submitted,



Robert L. Troike  
Attorney for Applicant  
Reg. No. 24,183

Texas Instruments Incorporated  
P.O. Box 655474, M/S 3999  
Dallas, TX 75265  
(202) 639-7710

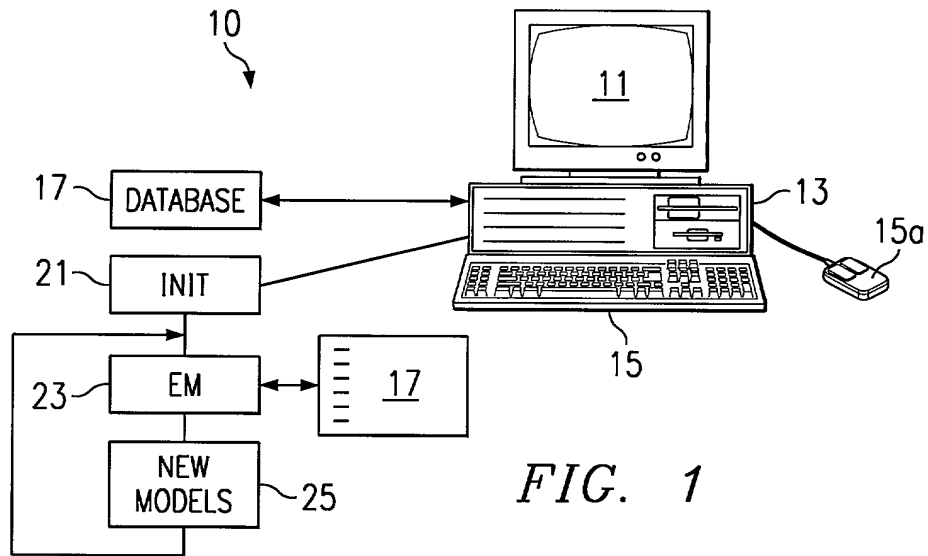


FIG. 1

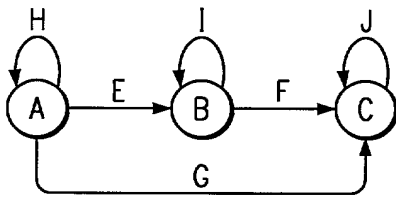


FIG. 2

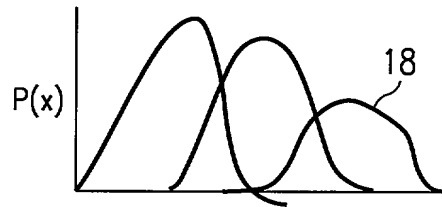


FIG. 3

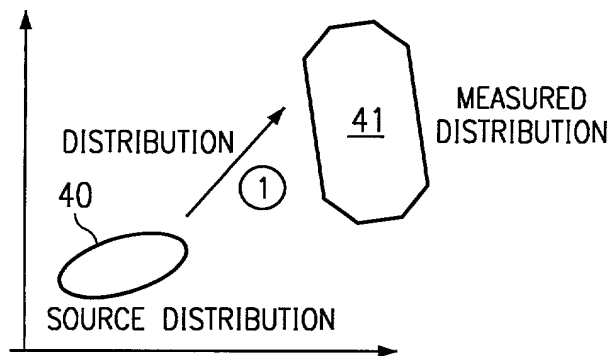


FIG. 4

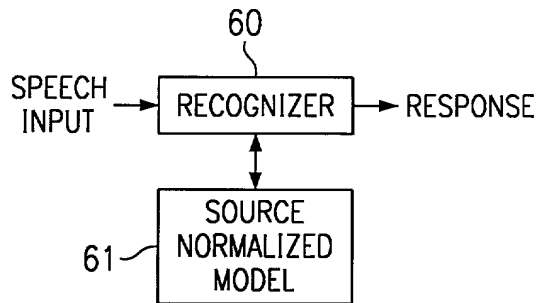


FIG. 6

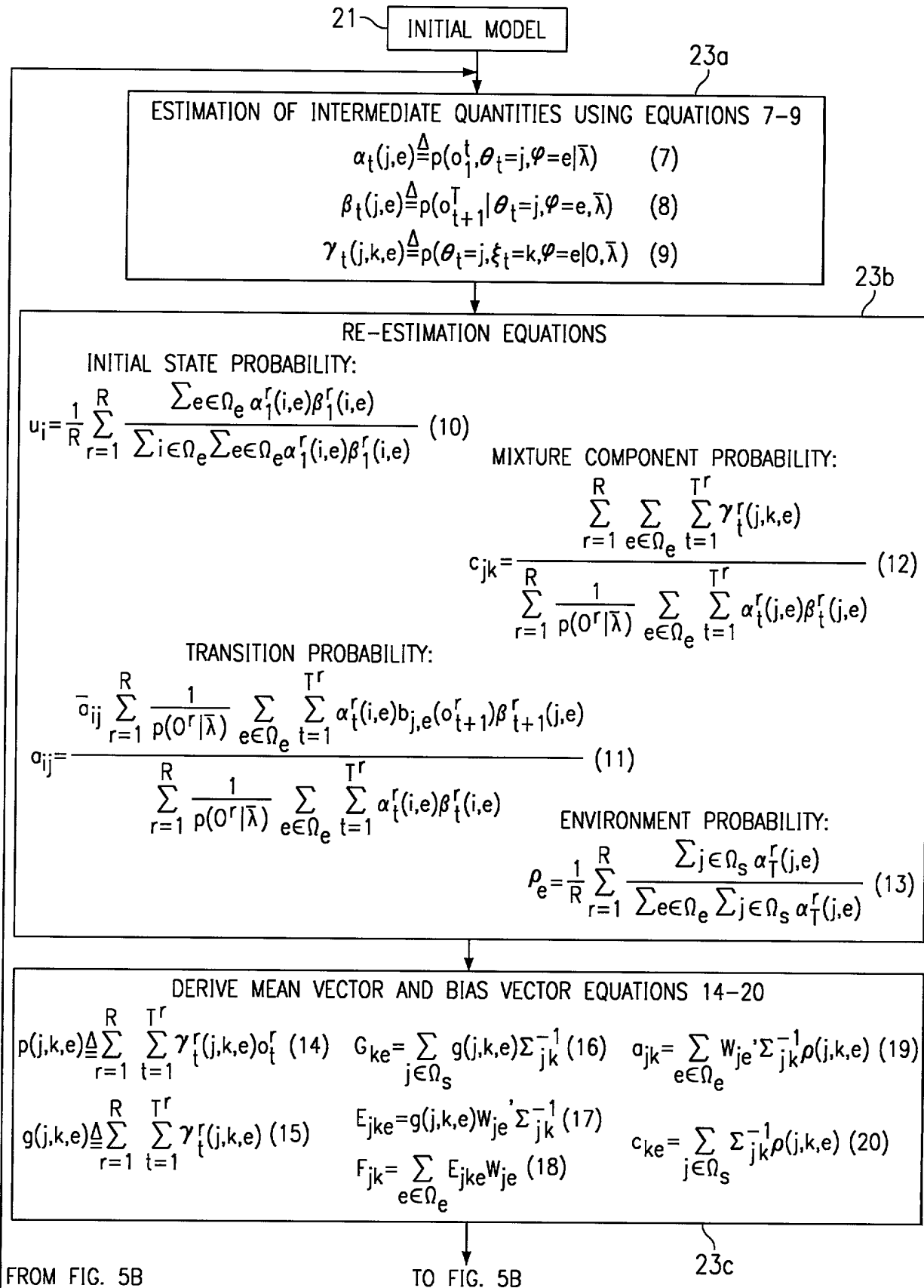


FIG. 5A

TO FIG. 5A

FROM FIG. 5A

23d

SOLVE JOINTLY FOR MEAN VECTORS AND BIAS VECTORS USING LINEAR EQUATIONS 21 AND 22

$$\sum_{e \in \Omega_e} E_{jke} b_{ke} + F_{jk} \mu_{jk} = a_{jk} \quad \forall j \in \Omega_e \quad (21)$$

$$G_{ke} b_{ke} + \sum_{j \in \Omega_s} H_{jke} \mu_{jk} = c_{ke} \quad \forall e \in \Omega_e \quad (22)$$

AND DETERMINE COVARIANCES USING EQUATION 23

$$\Sigma_{jk} = \frac{\sum_{e \in \Omega_e} \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j,k,e) \delta_t^r(j,k,e) \delta_t^r(j,e,k)}{\sum_{e \in \Omega_e} g(j,k,e)} \quad (23)$$

AND TRANSFORMATION USING EQUATIONS 24, 24, AND 26

$$Z_{je}^{(m)} = W_{je}^{(m)} R_{je}^{(m)} \quad (24)$$

$$Z_{je}^{(m,n)} \triangleq \sum_{k \in \Omega_m} \Sigma_{jk}^{-1(m,m)} \mu_{jk}^{(n)} \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j,k,e) (o_t^r - b_{ke})^{(m)} \quad (25)$$

$$R_{je}^{(p,n)}(m) \triangleq \sum_{k \in \Omega_m} \Sigma_{jk}^{-1(m,m)} \mu_{jk}^{(p)} \mu_{jk}^{(n)} \sum_{r=1}^R \sum_{t=1}^{T^r} \gamma_t^r(j,k,e) \quad (26)$$

REPLACE OLD MODEL PARAMETERS FOR THE CALCULATED ONES

24

YES

NEW MODEL

25

YES

CHANGE?

NO

27

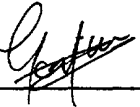
END

FIG. 5B

**APPLICATION FOR UNITED STATES PATENT**  
**DECLARATION AND POWER OF ATTORNEY**

As a below named inventor, I declare that my residence, post office address and citizenship are as stated below next to my name; that I verily believe that I am the original, first and sole inventor if only one name is listed below, or an original, first and joint inventor if plural inventors are named below, of the subject matter which is claimed and for which a patent is sought on the invention entitled as set forth below, and the title as set forth below which is described in the attached specification; that I have reviewed and understand the contents of the specification, including the claims, as amended by any amendment specifically referred to in the oath or declaration; that no application for patent or inventor's certificate on this invention has been filed by me or my legal representatives or assigns in any country foreign to the United States of America prior to the filing date of said application; and that I acknowledge my duty to disclose information which is material to the patentability of this application in accordance with Title 37, Code of Federal Regulations, section 1.56;

I further declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true, and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under section 1001 of Title 18 of the United States Code, and that such willful false statements may jeopardize the validity of the application or any patent issuing thereon.

TITLE OF INVENTION: Source Normalizaztion Training for HMM Modeling of Speech		
POWER OF ATTORNEY: I HEREBY APPOINT THE FOLLOWING ATTORNEYS TO PROSECUTE THIS APPLICATION AND TRANSACT ALL BUSINESS IN THE PATENT AND TRADEMARK OFFICE CONNECTED THEREWITH Robert L. Troike, #24,183; Richard L. Donaldson, #25,673; Jay M. Cantor, #19,906; Rene E. Grossman, #22,656, W. James Brady, III, #32,080; William B. Kempler, Reg. No. 28,228; Warren L. Franz, #28,716		
SEND CORRESPONDENCE TO: Robert L. Troike Texas Instruments Incorporated P.O. Box 655474, MS 219 Dallas, TX 75265		DIRECT TELEPHONE CALLS TO: Robert L. Troike 972/995-1364
NAME OF INVENTOR: (1) Yifan Gong	NAME OF INVENTOR: (2)	NAME OF INVENTOR: (3)
RESIDENCE & POST OFFICE ADDRESS: 7750 Walnut Hill Lane #2107 Dallas, Texas 75230	RESIDENCE & POST OFFICE ADDRESS:	RESIDENCE & POST OFFICE ADDRESS:
COUNTRY OF CITIZENSHIP: France	COUNTRY OF CITIZENSHIP:	COUNTRY OF CITIZENSHIP: -
SIGNATURE OF INVENTOR: 	SIGNATURE OF INVENTOR:	SIGNATURE OF INVENTOR:
DATE: September 29, 1997	DATE:	DATE: